



Full Bayesian inference with hazard mixture models

Julyan Arbel, Antonio Lijoi, Bernardo Nipoti

► To cite this version:

Julyan Arbel, Antonio Lijoi, Bernardo Nipoti. Full Bayesian inference with hazard mixture models. Computational Statistics and Data Analysis, 2016, 93, pp.359–372. 10.1016/j.csda.2014.12.003 . hal-01203296

HAL Id: hal-01203296

<https://hal.science/hal-01203296>

Submitted on 24 Sep 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Full Bayesian inference with hazard mixture models

Julyan Arbel^a, Antonio Lijoi^{b,a}, Bernardo Nipoti^{c,a,*}

^a*Collegio Carlo Alberto, via Real Collegio, 30, 10024 Moncalieri, Italy*

^b*Department of Economics and Management, University of Pavia, Via San Felice 5,
27100 Pavia, Italy*

^c*Department of Economics and Statistics, University of Torino, C.so Unione Sovietica
218/bis, 10134 Torino, Italy*

Abstract

Bayesian nonparametric inferential procedures based on Markov chain Monte Carlo marginal methods typically yield point estimates in the form of posterior expectations. Though very useful and easy to implement in a variety of statistical problems, these methods may suffer from some limitations if used to estimate non-linear functionals of the posterior distribution. The main goal is to develop a novel methodology that extends a well-established marginal procedure designed for hazard mixture models, in order to draw approximate inference on survival functions that is not limited to the posterior mean but includes, as remarkable examples, credible intervals and median survival time. The proposed approach relies on a characterization of the posterior moments that, in turn, is used to approximate the posterior distribution by means of a technique based on Jacobi polynomials. The inferential performance of this methodology is analysed by means of an extensive study of simulated data and real data consisting of leukemia remission times. Although tailored to the survival analysis context, the proposed procedure can be adapted to a range of other models for which moments of the posterior distribution can be estimated.

Keywords: Bayesian nonparametrics; Completely random measures; Hazard mixture models; Median survival time; Moment-based approximations; Survival analysis.

*Corresponding author at: Department of Economics and Statistics, University of Torino, C.so Unione Sovietica 218/bis, 10134 Torino, Italy. Phone: +39-011-6705023

Email addresses: julyan.arbel@carloalberto.org (Julyan Arbel), lijoi@unipv.it (Antonio Lijoi), bernardo.nipoti@carloalberto.org (Bernardo Nipoti)

1. Introduction

Most commonly used inferential procedures in Bayesian nonparametric practice rely on the implementation of sampling algorithms that can be gathered under the general umbrella of Blackwell–MacQueen Pólya urn schemes. These are characterized by the marginalization with respect to an infinite-dimensional random element that defines the de Finetti measure of an exchangeable sequence of observations or latent variables. Henceforth these will be referred to as *marginal methods*. Besides being useful for the identification of the basic building blocks of ready to use Markov chain Monte Carlo (MCMC) sampling strategies, marginal methods have proved to be effective for an approximate evaluation of Bayesian point estimators in the form of posterior means. They are typically used with models for which the predictive distribution is available in closed form. Popular examples are offered by mixtures of the Dirichlet process for density estimation (Escobar and West, 1995) and mixtures of gamma processes for hazard rate estimation (Ishwaran and James, 2004). While becoming well-established tools, these computational techniques are easily accessible also to practitioners through a straightforward software implementation (see for instance Jara et al., 2011). Though it is important to stress their relevance both in theory and in practice, it is also worth pointing out that Blackwell–MacQueen Pólya urn schemes suffer from some drawbacks which we wish to address here. Indeed, one easily notes that the posterior estimates provided by marginal methods are not suitably endowed with measures of uncertainty such as posterior credible intervals. Furthermore, using the posterior mean as an estimator is equivalent to choosing a square loss function whereas in many situations of interest other choices such as absolute error or 0–1 loss functions and, as corresponding estimators, median or mode of the posterior distribution of the survival function, at any fixed time point t , would be preferable. Finally, they do not naturally allow inference on functionals of the distribution of survival times, such as the median survival time, to be drawn. A nice discussion of these issues is provided by Gelfand and Kottas (2002) where the focus is on mixtures of the Dirichlet process: the authors suggest complementing the use of marginal methods with a sampling strategy that aims at generating approximate trajectories of the Dirichlet process from its truncated stick-breaking representation.

The aim is to propose a new procedure that combines closed-form analytical results arising from the application of marginal methods with an

approximation of the posterior distribution which makes use of posterior moments. The whole machinery is developed for the estimation of survival functions that are modeled in terms of hazard rate functions. To this end, let F denote the cumulative distribution function (CDF) associated to a probability distribution on \mathbb{R}^+ . The corresponding survival and cumulative hazard functions are denoted as

$$S(t) = 1 - F(t) \quad \text{and} \quad H(t) = - \int_{[0,t]} \frac{dF(s)}{F(s-)},$$

for any $t > 0$, respectively, where $F(s-) := \lim_{\varepsilon \downarrow 0} F(s - \varepsilon)$ for any positive s . If F is absolutely continuous, one has $H(t) = -\log(S(t))$ and the hazard rate function associated to F is, thus, defined as $h(t) = F'(t)/[1 - F(t-)]$. It should be recalled that survival analysis has been one of the most relevant areas of application of Bayesian nonparametric methodology soon after the groundbreaking contribution of [Ferguson \(1973\)](#). A number of papers in the '70s and the '80s have been devoted to the proposal of new classes of priors that accommodate for a rigorous analytical treatment of Bayesian inferential problems with censored survival data. Among these it is worth mentioning the neutral to the right processes proposed in [Doksum \(1974\)](#) and used to define a prior for the CDF F : since they share a conjugacy property they represent a tractable tool for drawing posterior inferences. Another noteworthy class of priors has been proposed in [Hjort \(1990\)](#), where a beta process is used as a nonparametric prior for the cumulative hazard function H has been proposed. Also in this case, one can considerably benefit from a useful conjugacy property.

As already mentioned, the plan consists in proposing a method for full Bayesian analysis of survival data by specifying a prior on the hazard rate h . The most popular example is the gamma process mixture that has been originally proposed in [Dykstra and Laud \(1981\)](#) and generalized in later work by [Lo and Weng \(1989\)](#) and [James \(2005\)](#) to include any mixing random measure and any mixed kernel. Recently [Lijoi and Nipoti \(2014\)](#) have extended such framework to the context of partially exchangeable observations. The uses of random hazard mixtures in practical applications have been boosted by the recent developments of powerful computational techniques that allow for an approximate evaluation of posterior inferences on quantities of statistical interest. Most of these arise from a marginalization with respect to a completely random measure that identifies the de Finetti measure of the exchangeable sequence of observations. See, e.g., [Ishwaran and James \(2004\)](#). Though they are quite simple to implement, the direct use of their output can only yield point estimation of the hazard rates, or

of the survival functions, at fixed time points through posterior means. The main goal of the present paper is to show that a clever use of a moment-based approximation method does provide a relevant upgrade on the type of inference one can draw via marginal sampling schemes. The takeaway message is that the information gathered by marginal methods is not confined to the posterior mean but is actually much richer and, if properly exploited, can lead to a more complete posterior inference. To understand this, one can refer to a sequence of exchangeable survival times $(X_i)_{i \geq 1}$ such that $\mathbb{P}[X_1 > t_1, \dots, X_n > t_n | \tilde{P}] = \prod_{i=1}^n \tilde{S}(t_i)$ where \tilde{P} is a random probability measure on \mathbb{R}^+ and $\tilde{S}(t) = \tilde{P}((t, \infty))$ is the corresponding random survival function. Given a suitable sequence of latent variables $(Y_i)_{i \geq 1}$, a closed-form expression for

$$\mathbb{E}[\tilde{S}^r(t) | \mathbf{X}, \mathbf{Y}], \quad \text{for any } r \geq 1, \text{ and } t > 0, \quad (1)$$

with $\mathbf{X} = (X_1, \dots, X_n)$ and $\mathbf{Y} = (Y_1, \dots, Y_n)$, will be provided. Our strategy consists in approximating the posterior distribution of $\tilde{S}(t)$, at each instant t , and relies on the fact that, along with the posterior mean, marginal models allow to straightforwardly estimate posterior moments of any order of $\tilde{S}(t)$. Indeed, an MCMC sampler yields a sample from the posterior distribution of \mathbf{Y} given \mathbf{X} : this can be used to integrate out the latent variables appearing in (1) and obtain a numerical approximate evaluation of the posterior moments $\mathbb{E}[\tilde{S}^r(t) | \mathbf{X}]$. These are finally used to deduce, with almost negligible effort, an approximation of the posterior distribution of $\tilde{S}(t)$ and, in turn, to estimate some meaningful functionals of $\tilde{S}(t)$.

It is to be mentioned that one could alternatively resort to a different approach that boils down to the simulation of the trajectories of the completely random measure that defines the underlying random probability measure from its posterior distribution. In density estimation problems, this is effectively illustrated in Nieto-Barajas et al. (2004), Nieto-Barajas and Prünster (2009) and Barrios et al. (2013). As for hazard rates mixtures estimation problems, one can refer to James (2005), Nieto-Barajas and Walker (2004) and Nieto-Barajas (2014). In particular, James (2005) provides a posterior characterization that is the key for devising a Ferguson and Klass (1972) representation of the posterior distribution of the completely random measure which enters the definition of the prior for the hazards. Some numerical aspects related to the implementation of the algorithm can be quite tricky since one needs to invert the Lévy intensity to simulate posterior jumps and a set of suitable latent variables need to be introduced in order to sample from the full conditionals of the hyperparameters. These aspects are well described and addressed in Nieto-Barajas (2014).

The paper is organized as follows. In Section 2 hazard mixture models are briefly reviewed together with some of their most important properties. Furthermore, explicit expressions characterizing the posterior moments of any order of a random survival function are provided both for general framework and for the extended gamma process case. Section 3 is dedicated to the problem of approximating the distribution of a random variable on $[0, 1]$, provided that the first N moments are known. In particular, a convenient methodology based on Jacobi polynomials is described in Section 3.1 and, then, implemented in Section 3.2 in order to approximate random survival functions. Its performance is tested through a thorough numerical investigation. The focus of Section 4 is the use of the introduced methodology for carrying out Bayesian inference on survival functions. Specifically, the algorithm is presented in Section 4.1 whereas simulated data and a real two-sample dataset on leukemia remission times are analysed in Sections 4.2 and 4.3 respectively. For the sake of exposition simplicity, technicalities such as expressions for the full conditional distributions involved in the algorithm and instructions on how to take into account the presence of censored data are postponed to the Appendix.

2. Hazard mixture models

A well-known nonparametric prior for the hazard rate function within multiplicative intensity models used in survival analysis arises as a mixture of *completely random measures* (CRMs). To this end, recall that a CRM $\tilde{\mu}$ on a space \mathbb{Y} is a boundedly finite random measure that, when evaluated at any collection of pairwise disjoint sets A_1, \dots, A_d , gives rise to mutually independent random variables $\tilde{\mu}(A_1), \dots, \tilde{\mu}(A_d)$, for any $d \geq 1$. Importantly, CRMs are almost surely discrete measures (Kingman, 1993). A detailed treatment on CRMs can also be found in Daley and Vere-Jones (2003). With reference to Theorem 1 in Kingman (1967), it is assumed that $\tilde{\mu}$ has no fixed atoms, which in turn implies the existence of a measure ν on $\mathbb{R}^+ \times \mathbb{Y}$ such that $\int_{\mathbb{R}^+ \times \mathbb{Y}} \min\{s, 1\} \nu(ds, dy) < \infty$ and

$$\mathbb{E} \left[e^{-\int_{\mathbb{Y}} f(y) \tilde{\mu}(dy)} \right] = \exp \left(- \int_{\mathbb{R}^+ \times \mathbb{Y}} [1 - \exp(-s f(y))] \nu(ds, dy) \right), \quad (2)$$

for any measurable function $f : \mathbb{Y} \rightarrow \mathbb{R}$ such that $\int_{\mathbb{Y}} |f| d\tilde{\mu} < \infty$, with probability 1. The measure ν is termed the *Lévy intensity* of $\tilde{\mu}$. For our purposes, it will be useful to rewrite ν as

$$\nu(ds, dy) = \rho_y(s) ds c P_0(dy),$$

where P_0 is a probability measure on \mathbb{Y} , c a positive parameter, and $\rho_y(s)$ is some transition kernel on $\mathbb{Y} \times \mathbb{R}^+$. If $\rho_y = \rho$, for any y in \mathbb{Y} , the CRM $\tilde{\mu}$ is said *homogeneous*. Henceforth, it is further assumed that P_0 is non-atomic. A well-known example corresponds to $\rho_y(s) = \rho(s) = e^{-s}/s$, for any y in \mathbb{Y} , which identifies a so-called *gamma CRM*. With such a choice of the Lévy intensity, it can be seen, from (2), that for any A such that $P_0(A) > 0$, the random variable $\tilde{\mu}(A)$ is gamma distributed, with shape parameter 1 and rate parameter $cP_0(A)$. If $k(\cdot; \cdot)$ is a transition kernel on $\mathbb{R}^+ \times \mathbb{Y}$, a prior for h is the distribution of the random hazard rate (RHR)

$$\tilde{h}(t) = \int_{\mathbb{Y}} k(t; y) \tilde{\mu}(dy), \quad (3)$$

where $\tilde{\mu}$ is a CRM on \mathbb{Y} . It is worth noting that, if $\lim_{t \rightarrow \infty} \int_0^t \tilde{h}(s) ds = \infty$ with probability 1, then one can adopt the following model

$$\begin{aligned} X_i | \tilde{P} &\stackrel{\text{iid}}{\sim} \tilde{P} \\ \tilde{P}((\cdot, \infty)) &\stackrel{\text{d}}{=} \exp\left(-\int_0^\cdot \tilde{h}(s) ds\right) \end{aligned} \quad (4)$$

for a sequence of (possibly censored) survival data $(X_i)_{i \geq 1}$. This means that \tilde{h} in (3) defines a random survival function $t \mapsto \tilde{S}(t) = \exp(-\int_0^t \tilde{h}(s) ds)$. In this setting, [Dykstra and Laud \(1981\)](#) characterize the posterior distribution of the so-called *extended gamma process*: this is obtained when $\tilde{\mu}$ is a gamma CRM and $k(t; y) = \mathbb{1}_{(0,t]}(y) \beta(y)$ for some positive right-continuous function $\beta : \mathbb{R}^+ \rightarrow \mathbb{R}^+$. The same kind of result is proved in [Lo and Weng \(1989\)](#) for *weighted gamma processes* corresponding to RHRs obtained when $\tilde{\mu}$ is still a gamma CRM and $k(\cdot; \cdot)$ is an arbitrary kernel. Finally, a posterior characterization has been derived in [James \(2005\)](#) for any CRM $\tilde{\mu}$ and kernel $k(\cdot; \cdot)$.

We shall quickly display such a characterization since it represents the basic result our construction relies on. For the ease of exposition we confine ourselves to the case where all the observations are exact, the extension to the case that includes right-censored data being straightforward and detailed in [James \(2005\)](#). See also [Appendix C](#). For an n -sample $\mathbf{X} = (X_1, \dots, X_n)$ of exact data, the likelihood function equals

$$\mathcal{L}(\tilde{\mu}; \mathbf{X}) = e^{-\int_{\mathbb{Y}} K_{\mathbf{X}}(y) \tilde{\mu}(dy)} \prod_{i=1}^n \int_{\mathbb{Y}} k(X_i; y) \tilde{\mu}(dy), \quad (5)$$

where $K_t(y) = \int_0^t k(s; y) ds$ and $K_{\mathbf{X}}(y) = \sum_{i=1}^n K_{X_i}(y)$. A useful augmentation suggests introducing latent random variables $\mathbf{Y} = (Y_1, \dots, Y_n)$ such that the joint distribution of $(\tilde{\mu}, \mathbf{X}, \mathbf{Y})$ coincides with

$$e^{-\int_{\mathbb{Y}} K_{\mathbf{X}}(y) \tilde{\mu}(dy)} \prod_{i=1}^n k(X_i; Y_i) \tilde{\mu}(dY_i) Q(d\tilde{\mu}), \quad (6)$$

where Q is the probability distribution of the completely random measure $\tilde{\mu}$, characterized by the Laplace transform functional in (2) (see for instance Daley and Vere-Jones, 2003). The almost sure discreteness of $\tilde{\mu}$ implies there might be ties among the Y_i 's with positive probability. Therefore, the distinct values among \mathbf{Y} are denoted as (Y_1^*, \dots, Y_k^*) , where $k \leq n$, and, for any $j = 1, \dots, k$, $C_j := \{l : Y_l = Y_j^*\}$ with $n_j = \#C_j$ as the cardinality of C_j . Thus, the joint distribution in (6) may be rewritten as

$$e^{-\int_{\mathbb{Y}} K_{\mathbf{X}}(y) \tilde{\mu}(dy)} \prod_{j=1}^k \tilde{\mu}(dY_j^*)^{n_j} \prod_{i \in C_j} k(X_i; Y_j^*) Q(d\tilde{\mu}). \quad (7)$$

We introduce, also, the density function

$$f(s \mid \kappa, \xi, y) \propto s^\kappa e^{-\xi s} \rho_y(s) \mathbb{1}_{\mathbb{R}^+}(s) \quad (8)$$

for any $\kappa \in \mathbb{N}$ and $\xi > 0$. The representation displayed in (7), combined with results concerning disintegrations of Poisson random measures, leads to prove the following

PROPOSITION 1 (James, 2005) *Let \tilde{h} be a RHR as defined in (3). The posterior distribution of \tilde{h} , given \mathbf{X} and \mathbf{Y} , coincides with the distribution of the random hazard*

$$\tilde{h}^* + \sum_{j=1}^k J_j k(\cdot; Y_j^*), \quad (9)$$

where $\tilde{h}^*(\cdot) = \int_{\mathbb{Y}} k(\cdot; y) \tilde{\mu}^*(dy)$ and $\tilde{\mu}^*$ is a CRM without fixed points of discontinuity whose Lévy intensity is

$$\nu^*(ds, dy) = e^{-sK_{\mathbf{X}}(y)} \rho_y(s) ds cP_0(dy).$$

The jumps J_1, \dots, J_k are mutually independent and independent of $\tilde{\mu}^*$. Moreover, for every $j = 1, \dots, k$, the distribution of the jump J_j has density function $f(\cdot \mid n_j, K_{\mathbf{X}}(Y_j^*), Y_j^*)$ with f defined in (8).

See [Lijoi et al. \(2008\)](#) for an alternative proof of this result. The posterior distribution of \tilde{h} displays a structure that is common to models based on CRMs, since it consists of the combination of two components: one without fixed discontinuities and the other with jumps at fixed points. In this case, the points at which jumps occur coincide with the distinct values of the latent variables Y_1^*, \dots, Y_k^* . Furthermore, the distribution of the jumps J_j depends on the respective locations Y_j^* .

Beside allowing us to gain insight on the posterior distribution of \tilde{h} , Proposition 1 is also very convenient for simulation purposes. See, e.g., [Ishwaran and James \(2004\)](#). Indeed, (9) allows obtaining an explicit expression for the posterior expected value of $\tilde{S}(t)$ (or, equivalently, of $\tilde{h}(t)$), for any $t > 0$, conditionally on the latent variables \mathbf{Y} . One can, thus, integrate out the vector of latent variables \mathbf{Y} , by means of a Gibbs type algorithm, in order to approximately evaluate the posterior mean of $\tilde{S}(t)$ (or $\tilde{h}(t)$). As pointed out in next section, a combination of Proposition 1 and of the same Gibbs sampler we have briefly introduced actually allows moments of $\tilde{S}(t)$, of any order, to be estimated. We will make use of the first N of these estimated moments to approximate, for each $t > 0$, the posterior distribution of $\tilde{S}(t)$ and therefore to have the tools for drawing meaningful Bayesian inference. The choice of a suitable value for N will be discussed in Section 3.2.

As pointed out in the Introduction, one can, in line of principle, combine Proposition 1 with the Ferguson and Klass representation to undertake an alternative approach that aims at simulating the trajectories from the posterior distribution of the survival function. This can be achieved by means of a Gibbs type algorithm that involves sampling $\tilde{\mu}^*$ and Y_j^* , for $j = 1, \dots, k$, from the corresponding full conditional distributions. Starting from the simulated trajectories one could then approximately evaluate all the posterior quantities of interest. The latter is an important feature of the method based on the Ferguson and Klass representation, that is shared only in part by our proposal. Indeed, extending the moment-based procedure to estimate functionals of $\tilde{S}(t)$, although achievable in many cases of interest, is not always straightforward. For instance, in order to carry out inference based on the posterior distribution of the random hazard rate $\tilde{h}(t)$, one should start with the estimation of the posterior moments of $\tilde{h}(t)$ and adapt accordingly the methodology which throughout the paper is developed for $\tilde{S}(t)$. An illustration, with an application to survival analysis, is provided in [Nieto-Barajas \(2014\)](#) and it appears that the approach, though achievable, may be difficult to implement. The main non-trivial issues one has to deal with are the inversion of the Lévy measure, needed to sample the jumps, and the sam-

pling from the full conditionals of the hyperparameters. The latter has been addressed by Nieto-Barajas (2014) through a clever augmentation scheme that relies on a suitable collection of latent variables. The approach based on the simulation of trajectories is an example of non-marginal, or *conditional*, method since it does not rely on the marginalization with respect to the mixing CRM $\tilde{\mu}$.

In the next sections, attention will be mainly devoted to marginal methods with the aim of showing that they allow for a full Bayesian inference, beyond the usual evaluation of posterior means. The required additional effort to accomplish this task is minimal and boils down to computing a finite number of posterior moments of $\tilde{S}(t)$, at a given t . An approximate evaluation of these moments can be determined by resorting to (9) which yields closed-form expressions for the posterior moments of the random variable $\tilde{S}(t)$, conditionally on both the data \mathbf{X} and the latent variables \mathbf{Y} .

PROPOSITION 2 *For every $t > 0$ and $r > 0$,*

$$\begin{aligned} \mathbb{E}[\tilde{S}^r(t) | \mathbf{X}, \mathbf{Y}] = & \exp \left\{ -c \int_{\mathbb{R}^+ \times \mathbb{Y}} \left(1 - e^{-rK_t(y)s} \right) e^{-K_{\mathbf{X}}(y)s} \rho(s) ds P_0(dy) \right\} \\ & \times \prod_{j=1}^k \frac{1}{B_j} \int_{\mathbb{R}^+} \exp \left\{ -s \left(rK_t(Y_j^*) + K_{\mathbf{X}}(Y_j^*) \right) \right\} s^{n_j} \rho(s) ds, \end{aligned}$$

where $B_j = \int_{\mathbb{R}^+} s^{n_j} \exp \left\{ -sK_{\mathbf{X}}(Y_j^*) \right\} \rho(s) ds$, for $j = 1, \dots, k$.

Although the techniques that will be described hold true for any specification of $\tilde{\mu}$ and kernel $k(\cdot; \cdot)$, the proposed illustration will focus on the extended gamma process case (Dykstra and Laud, 1981). More specifically, we consider a kernel $k(t; y) = \mathbb{1}_{(0, t]}(y)\beta$, with $\beta > 0$. This choice of kernel is known to be suitable for modeling monotone increasing hazard rates and to give rise to a class of random hazard functions with nice asymptotic properties (De Blasi et al., 2009). Moreover, without loss of generality, it is assumed that $X_1 > X_2 > \dots > X_n$. For notational convenience, one further sets $X_0 \equiv \infty$, $X_{n+1} \equiv 0$, $\xi_l \equiv \sum_{i=1}^l X_i$, for any $l \geq 1$, and $\xi_0 \equiv 0$. The next Corollary displays an expression for the conditional moments corresponding to this prior specification.

COROLLARY 1 *For every $t > 0$ and $r > 0$,*

$$\begin{aligned} \mathbb{E}[\tilde{S}^r(t) | \mathbf{X}, \mathbf{Y}] &= \prod_{i=0}^n \exp \left\{ -c \int_{X_{i+1} \wedge t}^{X_i \wedge t} \log \left(1 + r \frac{t - y}{\xi_i - iy + 1/\beta} \right) P_0(dy) \right\} \\ &\times \prod_{j=1}^k \left(1 + r \frac{(t - Y_j^*) \mathbb{1}_{[Y_j^*, \infty)}(t)}{\sum_{i=1}^n (X_i - Y_j^*) \mathbb{1}_{[Y_j^*, \infty)}(X_i) + 1/\beta} \right)^{-n_j}. \end{aligned} \quad (10)$$

By integrating out the vector of latent variables \mathbf{Y} in (10) one obtains an estimate of the posterior moments of $\tilde{S}(t)$. To this end one can resort to a Gibbs type algorithm whose steps will be described in Section 4.1.

3. Approximated inference via moments

3.1. Moment-based density approximation and sampling

Recovering a probability distribution from the explicit knowledge of its moments is a classical problem in probability and statistics that has received great attention in the literature. See, e.g., [Provost \(2005\)](#), references and motivating applications therein. Our specific interest in the problem is motivated by the goal of determining an approximation of the density function of a distribution supported on $[0, 1]$ that equals the posterior distribution of a random survival function evaluated at some instant t . This is a convenient case since, as the support is a bounded interval, all the moments exist and uniquely characterize the distribution, see [Rao \(1965\)](#). Moment-based methods for density functions' approximation can be essentially divided into two classes, namely methods that exploit orthogonal polynomial series ([Provost, 2005](#)) and maximum entropy methods ([Csiszár, 1975](#); [Mead and Papanicolaou, 1984](#)). Both these procedures rely on systems of equations that relate the moments of the distribution with the coefficients involved in the approximation. For our purposes the use of orthogonal polynomial series turns out to be more convenient since it ensures faster computations as it involves uniquely linear equations. This property is particularly important in our setting since the same approximation procedure needs to be implemented a large number of times in order to approximate the posterior distribution of a random survival function. Moreover, as discussed in [Epifani et al. \(2009\)](#), maximum entropy techniques can lead to numerical instability.

Specifically, we work with Jacobi polynomials, a broad class which includes, among others, Legendre and Chebyshev polynomials. They are well suited for the expansion of densities with compact support contrary to other polynomials like Laguerre and Hermite which can be preferred for densities with infinite or semi-infinite support (see [Provost, 2005](#)). While the classical Jacobi polynomials are defined on $[-1, 1]$, a suitable transformation

of such polynomials is considered so that their support coincides with $[0, 1]$ and therefore matches the support of the density we aim at approximating. That is, we consider a sequence of polynomials $(G_i)_{i \geq 0}$ such that, for every $i \in \mathbb{N}$, G_i is a polynomial of order i defined by $G_i(s) = \sum_{r=0}^i G_{i,r} s^r$, with $s \in [0, 1]$. The coefficients $G_{i,r}$ can be defined by a recurrence relation (see for example [Szegő, 1967](#)). Such polynomials are orthogonal with respect to the L^2 -product

$$\langle F, G \rangle = \int_0^1 F(s)G(s)w_{a,b}(s)ds,$$

where

$$w_{a,b}(s) = s^{a-1}(1-s)^{b-1}$$

is named *weight function* of the basis. Moreover, without loss of generality, the G_i 's can be assumed to be normalized and, therefore, $\langle G_i, G_j \rangle = \delta_{ij}$ for every $i, j \in \mathbb{N}$, where δ_{ij} is the Kronecker symbol. Any univariate density f supported on $[0, 1]$ can be uniquely decomposed on such a basis and therefore there is a unique sequence of real numbers $(\lambda_i)_{i \geq 0}$ such that

$$f(s) = w_{a,b}(s) \sum_{i=0}^{\infty} \lambda_i G_i(s). \quad (11)$$

Let us now consider a random variable S whose density f has support on $[0, 1]$. Its raw moments will be denoted by $\mu_r = \mathbb{E}[S^r]$, with $r \in \mathbb{N}$. From the evaluation of $\int_0^1 f(s) G_i(s) ds$ it follows that each λ_i coincides with a linear combination of the first i moments, specifically $\lambda_i = \sum_{r=0}^i G_{i,r} \mu_r$. Then, the polynomial approximation method consists in truncating the sum in (11) at a given level $i = N$. This procedure leads to a methodology that makes use only of the first N moments and provides the approximation

$$f_N(s) = w_{a,b}(s) \sum_{i=0}^N \left(\sum_{r=0}^i G_{i,r} \mu_r \right) G_i(s). \quad (12)$$

It is important to stress that the polynomial expansion approximation (12) is not necessarily a density as it might fail to be positive or to integrate to 1. In order to overcome this problem, the density π_N proportional to the positive part of f_N , i.e. $\pi_N(s) \propto \max(f_N(s), 0)$, will be considered. An importance sampling algorithm (see, e.g., [Robert and Casella, 2004](#)) will be used to sample from π_N . This is a method for drawing independent weighted samples (ϖ_ℓ, S_ℓ) from a distribution proportional to a given non-negative function, that exempts us from computing the normalizing constant. More precisely,

the method requires to pick a proposal distribution p whose support contains the support of π_N . A natural choice for p is the Beta distribution proportional to the weight function $w_{a,b}$. The weights are then defined by $\varpi_\ell \propto \max(f_N(S_\ell), 0)/p(S_\ell)$ such that they add up to 1.

An important issue related to any approximating method refers to the quantification of the approximating error. As for the described polynomial approach, the error can be assessed for large N by applying the asymptotic results in [Alexits and Földes \(1961\)](#). Specifically, the convergence $f_N(s) \rightarrow f(s)$ for $N \rightarrow \infty$, for all $s \in (0, 1)$, implies $\pi_N(s) \rightarrow f(s)$ for $N \rightarrow \infty$. Thus, if S_N denotes a random variable with distribution π_N , then the following convergence in distribution to the target random variable S holds:

$$S_N \xrightarrow{\mathcal{D}} S \text{ as } N \rightarrow \infty.$$

However, here the interest is in evaluating whether few moments allow for a good approximation of the posterior distribution of $\tilde{S}(t)$. This question will be addressed by means of an extensive numerical study in the next section. See [Epifani et al. \(2003\)](#) and [Epifani et al. \(2009\)](#) for a similar treatment referring to functionals of neutral-to-the-right priors and Dirichlet processes respectively.

3.2. Numerical study

In this section the quality of the approximation procedure described above is assessed by means of a simulation study. The rationale of the analysis consists in considering random survival functions for which moments of any order can be explicitly evaluated at any instant t , and then compare the true distribution with the approximation obtained by exploiting the knowledge of the first N moments. This in turn will provide an insight on the impact of N on the approximation error. To this end three examples of random survival functions will be considered, namely \tilde{S}_j with $j = 1, 2, 3$. For the illustrative purposes of this Section, it suffices to specify the distribution of the random variable that coincides with \tilde{S}_j evaluated in t , for every $t > 0$. Specifically, we consider a Beta, a mixture of Beta, and a normal distribution truncated to $[0, 1]$, that is

$$\begin{aligned}\tilde{S}_1(t) &\sim \text{Be}\left(\frac{S_0(t)}{a_1}, \frac{1 - S_0(t)}{a_1}\right), \\ \tilde{S}_2(t) &\sim \frac{1}{2}\text{Be}\left(\frac{S_0(t)}{a_2}, \frac{1 - S_0(t)}{a_2}\right) + \frac{1}{2}\text{Be}\left(\frac{S_0(t)}{a_3}, \frac{1 - S_0(t)}{a_3}\right), \\ \tilde{S}_3(t) &\sim \text{t}\mathcal{N}_{[0,1]}\left(S_0(t), \frac{S_0(t)(1 - S_0(t))}{a_4}\right),\end{aligned}$$

where $S_0(t) = e^{-t}$ and we have set $a_1 = 20$, $(a_2, a_3) = (10, 30)$ and $a_4 = 2$. Observe that, for every $t > 0$, $\mathbb{E}[\tilde{S}_1(t)] = \mathbb{E}[\tilde{S}_2(t)] = S_0(t)$ but the same does not hold true for $\tilde{S}_3(t)$.

For each $j = 1, 2, 3$, the first 10 moments of $\tilde{S}_j(t)$ were computed on a grid $\{t_1, \dots, t_{50}\}$ of 50 equidistant values of t in the range $[0, 2.5]$. The choice of working with 10 moments will be motivated at the end of the section. The importance sampler described in Section 3.1 was then used to obtain samples of size 10 000 from the distribution of $\tilde{S}_j(t_i)$, for each $j = 1, 2, 3$ and $i = 1, \dots, 50$. In Figure 1, for each \tilde{S}_j , we plot the true mean as well as the 95% highest density intervals for the true distribution and for the approximated distribution obtained by exploiting 10 moments. Notice that the focus is not on approximating the mean since moments of any order are the starting point of our procedure. Interestingly, the approximated intervals show a very good fit to the true ones in all the three examples. As for the Beta case, the fit is exact since the Beta-shaped weight function allows the true density to be recovered with the first two moments. As for the mixture of Beta, exact and approximated intervals can hardly be distinguished. Finally, the fit is pretty good also for the intervals in the truncated normal example. Similarly, in Figure 2 the true and the approximated densities of each $\tilde{S}_j(t)$ are compared for fixed t in $\{0.1, 0.5, 2.5\}$. Again, all the three examples show a very good pointwise fit.

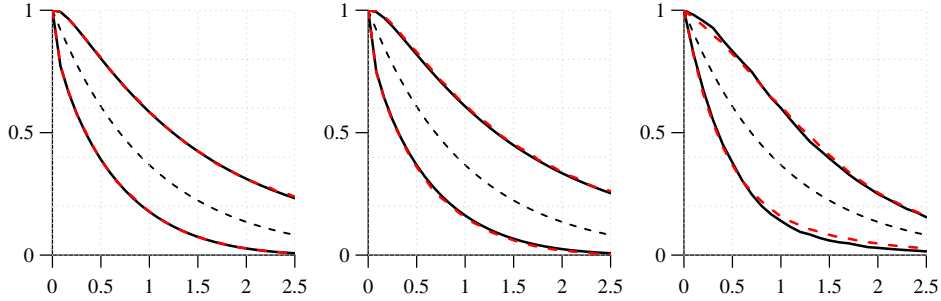


Figure 1: Mean of $\tilde{S}_j(t)$ (dashed black) and 95% highest density intervals for the true distribution (solid black) and the approximated distribution (dashed red) for the Beta ($j = 1$), mixture of Beta ($j = 2$) and truncated normal ($j = 3$) examples (left, middle and right, respectively).

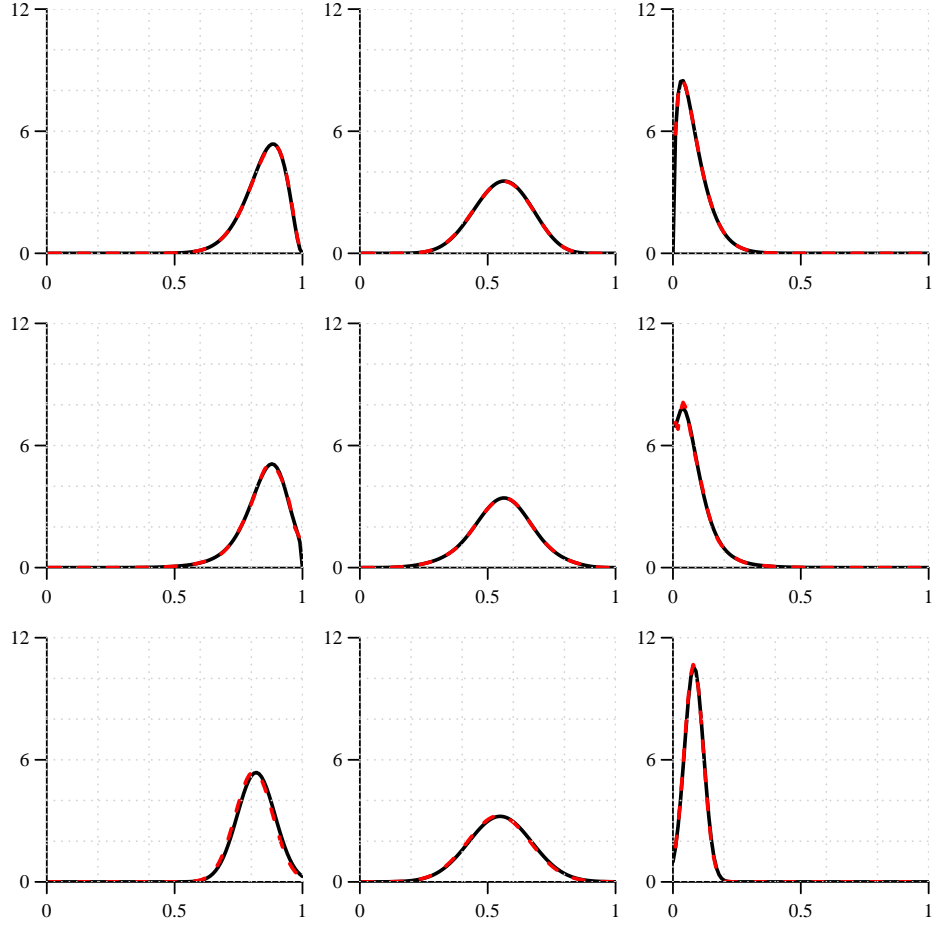


Figure 2: True density (solid black) and approximated one (dashed red) at time values $t = 0.1$ (left column), $t = 0.5$ (middle column) and $t = 2.5$ (right column), for the Beta ($j = 1$, top row), mixture of Beta ($j = 2$, middle row) and truncated normal ($j = 3$, bottom row) examples.

This section is concluded by assessing how the choice of N affects the approximation error. To this end, for each instant t on the grid, the true and approximated distributions of $\tilde{S}_j(t)$ are compared by computing the integrated squared error (L^2 error) between the two. Thus the average of these values is considered as a measure of the overall error of approximation. The results are illustrated in Figure 3. As expected, the approximation is exact in the Beta example. In the two other cases, it can be observed that the higher is the number of exploited moments, the lower is the average

approximation error. Nonetheless, it is apparent that the incremental gain of using more moments is more substantial when N is small whereas it is less impactful as N increases: for example in the mixture of Beta case, the L^2 error is 2.11, 0.97, 0.38 and 0.33 with N equal to 2, 4, 10 and 20 respectively. Moreover, when using a large number of moments, e.g. $N > 20$, some numerical instability can occur. These observations suggest that working with $N = 10$ moments in (12) strikes a good balance between accuracy of approximation and numerical stability.

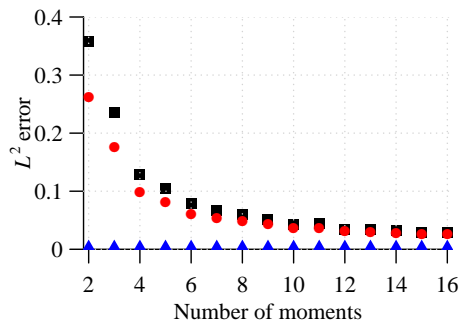


Figure 3: Average across t of the L^2 error between the true and the approximated densities of $\tilde{S}_j(t)$, in the Beta example (blue triangles), the mixture of Beta (red dots) and the truncated normal example (black squares). The approximation is exact in the Beta example.

4. Bayesian inference

In this section the characterization of the posterior moments of $\tilde{S}(t)$ provided in Proposition 2 is combined with the approximation procedure described in Section 3.1. The model specification (4) is completed by assuming an extended gamma prior for $\tilde{h}(t)$, with exponential base measure $P_0(dy) = \lambda \exp(-\lambda y) dy$, and considering the hyperparameters c and β random. This leads to the expression (A.1) for the posterior characterization of the moments. Finally we choose for both c and β independent gamma prior distributions with shape parameter 1 and rate parameter $1/3$ (so to ensure large prior variance) and set $\lambda = 1$. Given a sample of survival times $\mathbf{X} = \{X_1, \dots, X_n\}$, the first N moments of the posterior distribution of $\tilde{S}(t)$ are estimated for t on a grid of q equally-spaced points $\{t_1, \dots, t_q\}$ in an interval $[0, M]$. Such estimates are then exploited to approximate the posterior distribution of $\tilde{S}(t_i)$ for $i = 1, \dots, q$. This allows us to devise an

algorithm for carrying out full Bayesian inference on survival data. In the illustrations the focus will be on the estimation of the median survival time and, at any given t in the grid, of the posterior mean, posterior median, posterior mode and credibility intervals for $\tilde{S}(t)$. The same approach can be, in principle, used to estimate other functionals of interest.

4.1. Algorithm

The two main steps needed in order to draw samples from the posterior distribution of $\tilde{S}(t)$, for any $t \in \{t_1, \dots, t_q\}$, are summarized in Algorithm 1. First a Gibbs sampler is performed to marginalize the latent variables \mathbf{Y} and the hyperparameters (c, β) out of (A.1) and therefore, for every $i = 1, \dots, q$, an estimate for the posterior moments $\mathbb{E}[\tilde{S}^r(t_i)|\mathbf{X}]$, with $r = 1, \dots, N$, is obtained. The algorithm was run for $l_{\max} = 100\,000$ iterations, with a burn-in period of $l_{\min} = 10\,000$. Visual investigation of the traceplots of the parameters, in the illustrations of Sections 4.2 and 4.3, did not reveal any convergence issue. The second part consists in sampling from the posterior distribution of $\tilde{S}(t_i)$, for every $i = 1, \dots, q$, by means of the importance sampler described in Section 3.1. Specifically $\ell_{\max} = 10\,000$ values were sampled for each t_i on the grid.

The drawn samples allow us to approximately evaluate the posterior distribution of $\tilde{S}(t_i)$, for every $i = 1, \dots, q$. This, in turn, can be exploited to carry out meaningful Bayesian inference (Algorithm 2). As a remarkable example, we consider the median survival time, denoted by m . The identity for the cumulative distribution function of m

$$\mathbb{P}(m \leq t|\mathbf{X}) = \mathbb{P}(\tilde{S}(t) \leq 1/2|\mathbf{X})$$

allows us to evaluate the CDF of m at each time point t_i as $c_i = \mathbb{P}(\tilde{S}(t_i) \leq 1/2|\mathbf{X})$. Then, the median survival time m can be estimated by means of the following approximation:

$$\hat{m} = \mathbb{E}_{\mathbf{X}}[m] = \int_0^\infty \mathbb{P}[m > t|\mathbf{X}] dt \approx \frac{M}{q-1} \sum_{i=1}^q (1 - c_i) \quad (14)$$

where the subscript \mathbf{X} in $\mathbb{E}_{\mathbf{X}}[m]$ indicates that the integral is with respect to the distribution of $\tilde{S}(\cdot)$ conditional to \mathbf{X} . Equivalently,

$$\hat{m} \approx \sum_{i=1}^q t_i (c_{i+1} - c_i), \quad (15)$$

Algorithm 1 Posterior sampling

Part 1. Gibbs sampler

- 1: set $l = 0$ and admissible values for latent variables and hyperparameters, i.e. $\{Y_1 = Y_1^{(0)}, \dots, Y_n = Y_n^{(0)}\}$, $c = c^{(0)}$ and $\beta = \beta^{(0)}$
- 2: while $l < l_{\max}$, set $l = l + 1$, and
 - update $Y_j = Y_j^{(l)}$ by means of (B.1), for every $j = 1, \dots, n$
 - update $c = c^{(l)}$ and $\beta = \beta^{(l)}$ by means of (B.2) and (B.3)
 - if $l > l_{\min}$, compute

$$\mu_{r,t}^{(l)} = \mathbb{E}[\tilde{S}^r(t) | \mathbf{X}, \mathbf{Y}^{(l)}, c^{(l)}, \beta^{(l)}] \quad (13)$$

by means of (A.1) for each $r = 1, \dots, N$ and for each t in the grid

- 3: for each $r = 1, \dots, N$ and each t define $\hat{\mu}_{r,t} = \frac{1}{l_{\max} - l_{\min}} \sum_{l=l_{\min}+1}^{l_{\max}} \mu_{r,t}^{(l)}$

Part 2. Importance sampler

- 1: for each t , use (12) and define the approximate posterior density of $\tilde{S}(t)$ by $f_{N,t}(\cdot) = w_{a,b}(\cdot) \sum_{i=0}^N \left(\sum_{r=0}^i G_{i,r} \hat{\mu}_{r,t} \right) G_i(\cdot)$, where $\hat{\mu}_{0,t} \equiv 1$
 - 2: draw a weighted posterior sample $(\varpi_{\ell,t}, S_{\ell,t})_{\ell=1, \dots, \ell_{\max}}$ of $\tilde{S}(t)$, of size ℓ_{\max} , from $\pi_{N,t}(\cdot) \propto \max(f_{N,t}(\cdot), 0)$ by means of the important sampler described in Section 3.1
-

with the proviso that $c_{q+1} \equiv 1$. Moreover, the sequence $(c_i)_{i=1}^q$ can be used to devise credible intervals for the median survival time, cf. Part 1 of Algorithm 2. Note that both in (14) and in (15) the integrals on the left-hand-side are approximated by means of simple Riemann sums and the quality of such an approximation clearly depends on the choice of q and on M . Nonetheless, our investigations suggest that if q is sufficiently large the estimates we obtain are pretty stable and that the choice of M is not crucial since, for t_i sufficiently large, the term $1 - c_i$ involved in (14) is approximately equal to 0. Finally, the posterior samples generated by Algorithm 1 can be used to obtain a t -by- t estimation of other functionals that convey meaningful information such as the posterior mode and median (together with the posterior mean), cf. Part 2 of Algorithm 2.

Algorithm 2 Bayesian inference

Part 1. Median survival time

- 1: use the weighted sample $(\varpi_{\ell,t_i}, S_{\ell,t_i})_{\ell=1,\dots,\ell_{\max}}$ to estimate, for each $i = 1, \dots, q$, $c_i = \mathbb{P}(\tilde{S}(t_i) \leq 1/2 | \mathbf{X})$
- 2: plug the c_i 's in (15) to obtain \hat{m}
- 3: use the sequence $(c_i)_{i=1}^q$ as a proxy for the posterior distribution of m so to devise credible intervals for \hat{m} .

Part 2. t -by- t functionals

- 1: use the weighted sample $(\varpi_{\ell,t_i}, S_{\ell,t_i})_{\ell=1,\dots,\ell_{\max}}$ to estimate, for each $i = 1, \dots, q$, $a_i = \inf_{x \in [0,1]} \{\mathbb{P}(\tilde{S}(t_i) \leq x | \mathbf{X}) \geq 1/2\}$ and $b_i = \text{mode}\{\tilde{S}(t_i) | \mathbf{X}\}$
 - 2: use the sequences $(a_i)_{i=1}^q$ and $(b_i)_{i=1}^q$ to approximately evaluate, t -by- t , posterior median and mode respectively
 - 3: use the weighted sample $(\varpi_{\ell,t_i}, S_{\ell,t_i})_{\ell=1,\dots,\ell_{\max}}$ to devise t -by- t credible intervals
-

The rest of this section is divided in two parts in which Algorithms 1 and 2 are applied to analyse simulated and real survival data. In Section 4.2 the focus is on the estimation of the median survival time for simulated samples of varying size. In Section 4.3 we analyse a real two-sample dataset and we estimate posterior median and mode, together with credible intervals, of $\tilde{S}(t)$. In both illustrations our approximations are based on the first $N = 10$ moments.

4.2. Application to simulated survival data

Consider four samples of size $n = 25, 50, 100, 200$, from a mixture f of Weibull distributions, defined by

$$f = \frac{1}{2} \text{Wbl}(2, 2) + \frac{1}{2} \text{Wbl}(2, 1/2).$$

After observing that the largest observation in the samples is 4.21, we set $M = 5$ and $q = 100$ for the analysis of each sample. By applying Algorithms 1 and 2 we approximately evaluate, t -by- t , the posterior distribution of $\tilde{S}(t)$ together with the posterior distribution of the median survival time m . In Figure 4 the focus is on the sample corresponding to $n = 100$. On the left panel, true survival function and Kaplan–Meier estimate are plotted. By investigating the right panel it can be appreciated that the estimated

HPD credible regions for $\tilde{S}(t)$ contain the true survival function. Moreover, the posterior distribution of m is nicely concentrated around the true value $m_0 = 0.724$.

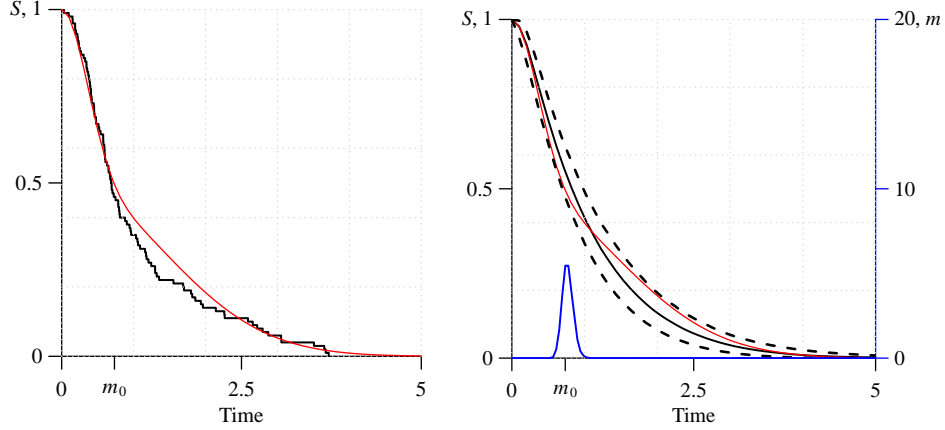


Figure 4: (Simulated dataset, $n = 100$.) Left: true survival function (red line) and Kaplan–Meier estimate (black line). Right: true survival function (red line) and estimated posterior mean (black solid line) with 95% HPD credible intervals for $\tilde{S}(t)$ (black dashed lines); the blue plot appearing in the panel on the right is the posterior distribution of the median survival time m .

The performance of the introduced methodology is investigated as the sample size n grows. Table 1 summarizes the values obtained for \hat{m} and the corresponding credible intervals. For all the sample sizes considered, credible intervals for \hat{m} contain the true value. Moreover, as expected, as n grows, they shrink around m_0 : for example the length of the interval reduces from 0.526 to 0.227 when the size n changes from 25 to 200. Finally, for all these samples, the estimated median survival time \hat{m} is closer to m_0 than the empirical estimator \hat{m}_e .

4.3. Application to real survival data

The described methodology is now used to analyse a well known two-sample dataset involving leukemia remission times, in weeks, for two groups of patients, under active drug treatment and placebo respectively. The same dataset was studied, e.g., by Cox (1972). Observed remission times for patients under treatment (T) are

$\{6, 6, 6, 6^*, 7, 9^*, 10, 10^*, 11, 13, 16, 17^*, 19^*, 20^*, 22, 23, 25^*, 32^*, 32^*, 34^*, 35^*\}$,

Table 1: (Simulated datasets.) Comparison of the estimated median survival time (\hat{m}) obtained by means of our Bayesian nonparametric procedure (BNP) and the empirical median survival time \hat{m}_e , for different sample sizes. For BNP estimation we show \hat{m} , the absolute error $|\hat{m} - m_0|$ and the 95%-credible interval (CI); last two columns show the empirical estimate \hat{m}_e and the corresponding absolute error $|\hat{m}_e - m_0|$. The true median survival time m_0 is 0.724.

sample size	BNP			Empirical	
	\hat{m}	error	CI	\hat{m}_e	error
25	0.803	0.079	(0.598, 1.124)	0.578	0.146
50	0.734	0.010	(0.577, 0.967)	0.605	0.119
100	0.750	0.026	(0.622, 0.912)	0.690	0.034
200	0.746	0.022	(0.669, 0.896)	0.701	0.023

where stars denote right-censored observations. Details on the censoring mechanism and on how to adapt the methodology to right-censored observations are provided in [Appendix C](#). On the other side, remission times of patients under placebo (P) are all exact and coincide with

$$\{1, 1, 2, 2, 3, 4, 4, 5, 5, 8, 8, 8, 11, 11, 12, 12, 15, 17, 22, 23\}.$$

For this illustration we set $M = 2 \max(\mathbf{X})$, that is $M = 70$, and $q = 50$. For both samples posterior mean, median and mode as well as 95% credible intervals, are estimated and compared. In the left panel of [Figure 5](#) such estimates are plotted for sample T. By inspecting the plot, it is apparent that, for large values of t , posterior mean, median and mode show significantly different behaviors, with posterior mean being more optimistic than posterior median and mode. It is worth stressing that such differences, while very meaningful for clinicians, could not be captured by marginal methods for which only the posterior mean would be available. A fair analysis must take into account the fact that, up to $t = 23$, i.e. the value corresponding to the largest non-censored observation, the three curves are hardly distinguishable. The different patterns for larger t might therefore depend on the prior specification of the model. Nonetheless, this example is meaningful as it shows that a more complete posterior analysis is able to capture differences, if any, between posterior mean, median and mode.

When relying on marginal methods, the most natural choice for estimating the uncertainty of posterior estimates consists in considering the quantiles intervals corresponding to the output of the Gibbs sampler, that we refer to as *marginal intervals*. This leads to consider, for any fixed t , the interval

whose lower and upper extremes are the quantiles of order 0.025 and 0.975, respectively, of the sample of conditional moments $\{\mu_{1,t}^{(l_{\min}+1)}, \dots, \mu_{1,t}^{(l_{\max})}\}$ defined in (13). In the middle panel of Figure 5 the estimated 95% HPD intervals for $\hat{S}(t)$ and the marginal intervals corresponding to the output of the Gibbs sampler are compared. In this example, the marginal method clearly underestimates the uncertainty associated to the posterior estimates. This can be explained by observing that, since the underlying completely random measure has already been marginalized out, the intervals arising from the Gibbs sampler output, capture only the variability of the posterior mean that can be traced back to the latent variables \mathbf{Y} and the parameters (c, β) . As a result, the uncertainty detected by the marginal method leads to credible intervals that can be significantly narrower than the actual posterior credible intervals that we approximate through the moment-based approach. This suggests that the use of intervals produced by marginal methods as proxies for posterior credible intervals should be, in general, avoided.

The analysis is concluded by observing that the availability of credible intervals for survival functions can be of great help in comparing treatments. In the right panel of Figure 5 posterior means as well as corresponding 95% HPD intervals are plotted for both samples T and P. By inspecting the plot, for example, the effectiveness of the treatment seems clearly significant as, essentially, there is no overlap between credible intervals of the two groups.

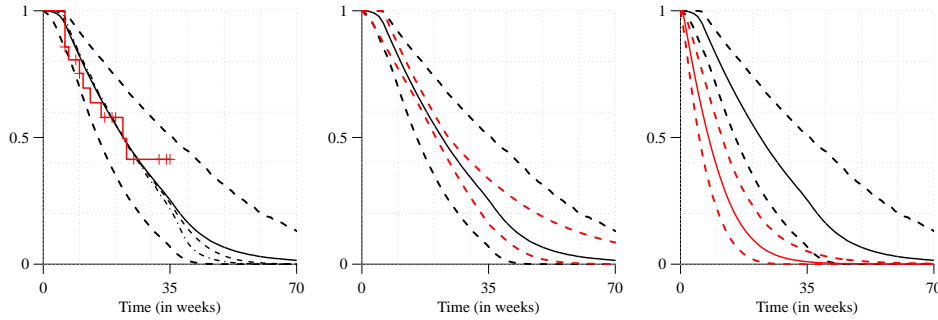


Figure 5: Left: comparison of posterior mean (solid line), median (dashed line) and mode (point dashed line) in dataset T, with 95% HPD credible intervals (dashed line). The Kaplan–Meier estimate is plotted in red. Middle: comparison of the 95% HPD credible interval (dashed black line) with the marginal interval (dashed red line). Right: comparison of samples T (black) and P (red), with posterior means (solid) and 95% HPD credible intervals (dashed).

Acknowledgment

J. Arbel and A. Lijoi are supported by the European Research Council (ERC) through StG “N-BNP” 306406.

Appendix A. Moments under exponential P_0

An explicit expression for (10) is provided when $P_0(dy) = \lambda \exp(-\lambda y)dy$ and the hyperparameters c and β are considered random.

$$\begin{aligned}
\mathbb{E}[\tilde{S}^r(t) | \mathbf{X}, \mathbf{Y}, c, \beta] = & \exp \left\{ -ce^{-f_{0,r}(0)} [\text{Ei}(f_{0,r}(t)) - \text{Ei}(f_{0,r}(X_1 \wedge t))] \right\} \left(\frac{f_{0,r}(X_1 \wedge t)}{f_{0,r}(t)} \right)^{-ce^{-\lambda(X_1 \wedge t)}} \\
& \times \prod_{i=1}^n \exp \left\{ -ce^{-f_{i,r}(0)} [\text{Ei}(f_{i,r}(X_i \wedge t)) - \text{Ei}(f_{i,r}(X_{i+1} \wedge t))] \right. \\
& \quad \left. -ce^{-f_{i,0}(0)} [\text{Ei}(f_{i,0}(X_{i+1} \wedge t)) - \text{Ei}(f_{i,0}(X_i \wedge t))] \right\} \\
& \times \left(\frac{i}{i+r} \frac{f_{i,0}(X_i \wedge t)}{f_{i,r}(X_i \wedge t)} \right)^{-ce^{-\lambda(X_i \wedge t)}} \left(\frac{i+r}{i} \frac{f_{i,r}(X_{i+1} \wedge t)}{f_{i,0}(X_{i+1} \wedge t)} \right)^{-ce^{-\lambda(X_{i+1} \wedge t)}} \\
& \times \prod_{j=1}^k \left(1 + r \frac{(t - Y_j^*) \mathbb{1}_{[Y_j^*, \infty)}(t)}{\sum_{i=1}^n (X_i - Y_j^*) \mathbb{1}_{[Y_j^*, \infty)}(X_i) + 1/\beta} \right)^{-n_j}, \quad (\text{A.1})
\end{aligned}$$

where $\text{Ei}(\cdot)$ is the exponential integral function defined for non-zero real values z by

$$\text{Ei}(z) = - \int_{-z}^{\infty} \frac{e^{-t}}{t} dt$$

and the function $f_{i,r}$, for $i, r \geq 0$ such that $i + r > 0$, is defined by

$$f_{i,r}(x) = \lambda \left(\frac{\xi_i + 1/\beta + rt}{i+r} - x \right).$$

Appendix B. Full conditional distributions

In this section we provide expressions for the full conditional distributions needed in the algorithm described in Section 4.1 for extended gamma processes with base measure $P_0(dy) = \lambda \exp(-\lambda y)dy$. These distributions are easily derived, up to a constant, from the joint distribution of the vector

$(\mathbf{X}, \mathbf{Y}, c, \beta)$, that can be obtained from (7). Therefore we start by providing the full conditional distribution for the latent variable Y_i , with $i = 1, \dots, n$, where $\mathbf{Y}^{(-i)}$ denotes the vector of distinct values $(\tilde{Y}_1^*, \dots, \tilde{Y}_{k^*}^*)$ in $(Y_1, \dots, Y_{i-1}, Y_{i+1}, \dots, Y_n)$ and $(n_1^{(-i)}, \dots, n_{k^*}^{(-i)})$ represent the corresponding frequencies.

$$\mathbb{P}[Y_i = dy \mid \mathbf{X}, \mathbf{Y}^{(-i)}, c, \beta] = p_0 G_0(dy) + \sum_{j=1}^{k^*} p_j \delta_{\tilde{Y}_j^*}(dy), \quad (\text{B.1})$$

where

$$p_0 \propto c \lambda \sum_{j=i}^n \frac{1}{j} e^{-\lambda \frac{\xi_j + 1/\beta}{j}} [\text{Ei}(f_{j,0}(X_{j+1})) - \text{Ei}(f_{j,0}(X_j))],$$

$$p_j \propto \mathbb{1}_{\{Y_j^* \leq X_i\}} \frac{n_j^{(-i)}}{\sum_{l=1}^n (X_l - \tilde{Y}_j^*) \mathbb{1}_{[0, X_l)}(\tilde{Y}_j^*) + 1/\beta}$$

and

$$G_0(dy) \propto \mathbb{1}_{[0, X_i)}(y) e^{-\lambda y} \frac{1}{\sum_{j=1}^n (X_j - y) \mathbb{1}_{[0, X_j)}(y) + 1/\beta} dy.$$

Finally, the full conditional distributions for the parameters c and β are given respectively by

$$\begin{aligned} \mathcal{L}(c \mid \mathbf{X}, \mathbf{Y}, \beta) &\propto \mathcal{L}_0(c) c^k \beta^{-c} \prod_{i=1}^n \exp \left\{ -ce^{-f_{i,0}(0)} [\text{Ei}(f_{i,0}(X_i)) - \text{Ei}(f_{i,0}(X_{i+1}))] \right\} \\ &\quad \times \frac{(\xi_i + 1/\beta - iX_{i+1})^{-ce^{-\lambda X_{i+1}}}}{(\xi_i + 1/\beta - iX_i)^{-ce^{-\lambda X_i}}} \end{aligned} \quad (\text{B.2})$$

and

$$\begin{aligned} \mathcal{L}(\beta \mid \mathbf{X}, \mathbf{Y}, c) &\propto \mathcal{L}_0(\beta) \beta^{-c} \prod_{i=1}^n \exp \left\{ -ce^{-f_{i,0}(0)} [\text{Ei}(f_{i,0}(X_i)) - \text{Ei}(f_{i,0}(X_{i+1}))] \right\} \\ &\quad \times \frac{(\xi_i + 1/\beta - iX_{i+1})^{-ce^{-\lambda X_{i+1}}}}{(\xi_i + 1/\beta - iX_i)^{-ce^{-\lambda X_i}}} \prod_{j=1}^k \left(\sum_{i=1}^n (X_i - Y_j^*) \mathbb{1}_{[Y_j^*, \infty)}(X_i) + 1/\beta \right)^{-n_j}, \end{aligned} \quad (\text{B.3})$$

where $\mathcal{L}_0(c)$ and $\mathcal{L}_0(\beta)$ are the prior distributions of c and β respectively.

Appendix C. Censored observations

The methodology presented in Section 4 needs to be adapted to the presence of right-censored observations in order to be applied to the dataset in Section 4.3. Here we introduce some notation and illustrate how the posterior characterization of Proposition 1 changes when data are censored. To this end, let C_i be the right-censoring time corresponding to X_i , and define $\Delta_i = \mathbb{1}_{(0, C_i]}(X_i)$, so that Δ_i is either 0 or 1 according as to whether X_i is censored or exact. The actual i th observation is $T_i = \min(X_i, C_i)$ and, therefore, data consist of pairs $\mathbf{D} = \{(T_i, \Delta_i)\}_{i=1 \dots n}$. In this setting, the likelihood in (5) can be rewritten as

$$\mathcal{L}(\tilde{\mu}; \mathbf{D}) = e^{-\int_{\mathbb{Y}} K_{\mathbf{D}}^*(y) \tilde{\mu}(dy)} \prod_{i: \Delta_i=1} \int_{\mathbb{Y}} k(T_i; y) \tilde{\mu}(dy),$$

where

$$K_{\mathbf{D}}^*(y) = \sum_{i=1}^n \int_0^{T_i} k(s; y) ds.$$

By observing that the censored times are involved only through $K_{\mathbf{D}}^*$, the results derived in Proposition 1 under the assumption of exact data easily carry over to the case with right-censored data. The only changes refer to $K_{\mathbf{X}}$, that is replaced by $K_{\mathbf{D}}^*$, and the jump components which occur only at the distinct values of the latent variables that correspond to exact observations. For instance in Proposition 1, the Lévy intensity of the part of the CRM without fixed points of discontinuity is modified by

$$\nu^*(ds, dy) = e^{-sK_{\mathbf{D}}^*(y)} \rho_y(s) ds cP_0(dy),$$

while the distribution of the jump J_j has density function $f(\cdot | n_j^*, K_{\mathbf{D}}^*(Y_j^*), Y_j^*)$ with f defined in (8) and $n_j^* = \#\{i : Y_i = Y_j^* \text{ and } \Delta_i = 1\}$. Adapting the results of Proposition 2 and Corollary 1, as well as the full conditional distributions in Appendix B, is then straightforward.

References

- Alexits, G. and Földes, I. (1961). *Convergence problems of orthogonal series*. Pergamon Press New York.
- Barrios, E., Lijoi, A., Nieto-Barajas, L. E., Prünster, I., et al. (2013). Modeling with normalized random measure mixture models. *Statist. Sci.*, 28(3):313–334.

- Cox, D. (1972). Regression models and life tables (with discussion). *J. Roy. Stat. Soc. B Met.*, 34(2):187–202.
- Csiszár, I. (1975). I-divergence geometry of probability distributions and minimization problems. *Ann. Probab.*, 3(1):146–158.
- Daley, D. and Vere-Jones, D. (2003). *An introduction to the theory of point processes. Vol. I*. Springer-Verlag, New York.
- De Blasi, P., Peccati, G., Prünster, I., et al. (2009). Asymptotics for posterior hazards. *Ann. Statist.*, 37(4):1906–1945.
- Doksum, K. (1974). Tailfree and neutral random probabilities and their posterior distributions. *Ann. Probab.*, 2(2):183–201.
- Dykstra, R. and Laud, P. (1981). A Bayesian nonparametric approach to reliability. *Ann. Statist.*, 9(2):356–367.
- Epifani, I., Guglielmi, A., and Melilli, E. (2009). Moment-based approximations for the law of functionals of Dirichlet processes. *Applied Mathematical Sciences*, 3(20):979–1004.
- Epifani, I., Lijoi, A., and Prünster, I. (2003). Exponential functionals and means of neutral-to-the-right priors. *Biometrika*, 90(4):791–808.
- Escobar, M. D. and West, M. (1995). Bayesian density estimation and inference using mixtures. *J. Am. Stat. Assoc.*, 90(430):577–588.
- Ferguson, T. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.*, 1(2):209–230.
- Ferguson, T. S. and Klass, M. J. (1972). A representation of independent increment processes without gaussian components. *Ann. Math. Stat.*, 43(5):1634–1643.
- Gelfand, A. E. and Kottas, A. (2002). A computational approach for full nonparametric Bayesian inference under Dirichlet process mixture models. *J. Comput. Graph. Stat.*, 11(2):289–305.
- Hjort, N. (1990). Nonparametric Bayes estimators based on beta processes in models for life history data. *Ann. Statist.*, 18(3):1259–1294.
- Ishwaran, H. and James, L. (2004). Computational methods for multiplicative intensity models using weighted gamma processes: proportional hazards, marked point processes, and panel count data. *J. Am. Stat. Assoc.*, 99(465):175–190.

- James, L. (2005). Bayesian Poisson process partition calculus with an application to Bayesian Lévy moving averages. *Ann. Statist.*, 33(4):1771–1799.
- Jara, A., Hanson, T., Quintana, F., Müller, P., and Rosner, G. (2011). DPpackage: Bayesian non-and semi-parametric modelling in R. *J. Stat. Softw.*, 40(5):1.
- Kingman, J. F. C. (1967). Completely random measures. *Pacific J. Math.*, 21(1):59–78.
- Kingman, J. F. C. (1993). *Poisson processes*, volume 3. Oxford university press.
- Lijoi, A. and Nipoti, B. (2014). A class of hazard rate mixtures for combining survival data from different experiments. *J. Am. Stat. Assoc.*, 109(506):802–814.
- Lijoi, A., Prünster, I., and Walker, S. G. (2008). Posterior analysis for some classes of nonparametric models. *J. Nonparametr. Stat.*, 20(5):447–457.
- Lo, A. and Weng, C. (1989). On a class of Bayesian nonparametric estimates. II. Hazard rate estimates. *Ann. I. Stat. Math.*, 41(2):227–245.
- Mead, L. R. and Papanicolaou, N. (1984). Maximum entropy in the problem of moments. *J. Math. Phys.*, 25(8):2404–2417.
- Nieto-Barajas, L. E. (2014). Bayesian semiparametric analysis of short- and long-term hazard ratios with covariates. *Comput. Stat. Data. An.*, 71:477–490.
- Nieto-Barajas, L. E. and Prünster, I. (2009). A sensitivity analysis for Bayesian nonparametric density estimators. *Stat. Sinica*, 19:685–705.
- Nieto-Barajas, L. E., Prünster, I., and Walker, S. G. (2004). Normalized random measures driven by Increasing Additive Processes. *Ann. Statist.*, 32(6):2343–2360.
- Nieto-Barajas, L. E. and Walker, S. G. (2004). Bayesian nonparametric survival analysis driven by Lévy driven Markov processes. *Stat. Sinica*, 14:1127–1146.
- Provost, S. B. (2005). Moment-based density approximants. *Mathematica J.*, 9(4):727–756.

- Rao, C. R. (1965). *Linear statistical inference and its applications*. Wiley, New York.
- Robert, C. P. and Casella, G. (2004). *Monte Carlo statistical methods*. Springer-Verlag, New York.
- Szegő, G. (1967). *Orthogonal polynomials*. American Mathematical Society Colloquium Publications.